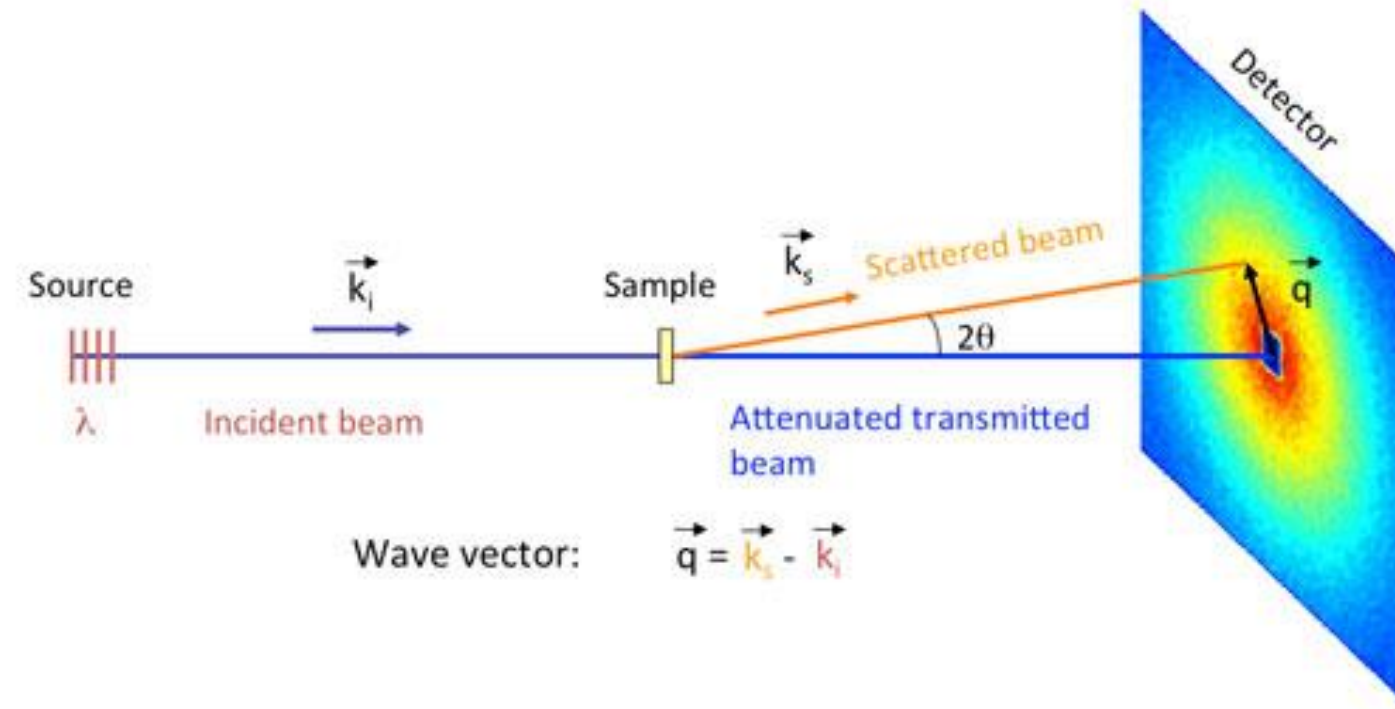# GPU PROGRAMMING: 2D FITTING FOR SMALL-ANGLE NEUTRON SCATTERING DATA
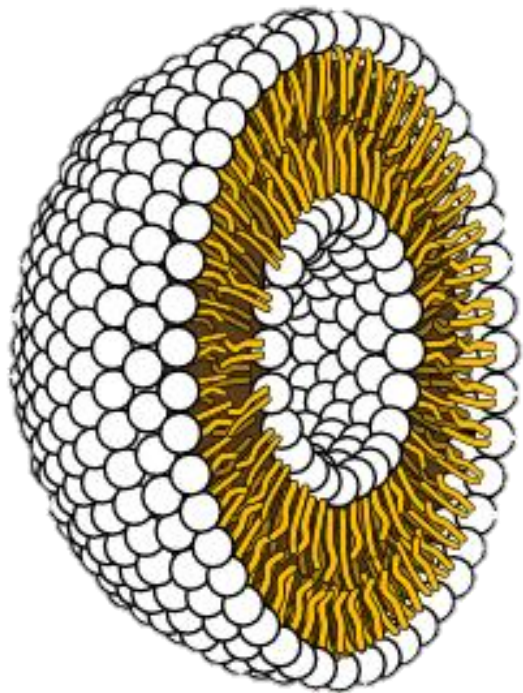
Helen Park

Paul Kienzle

Matt Wasbrough

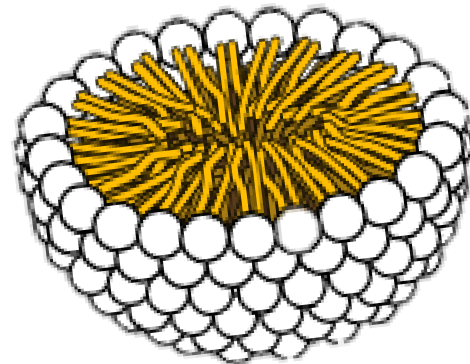# DATA USED: COLLECTED BY NEUTRON SCATTERING



- Neutrons scatter off large scale structures in the material
- Reveal structure of material
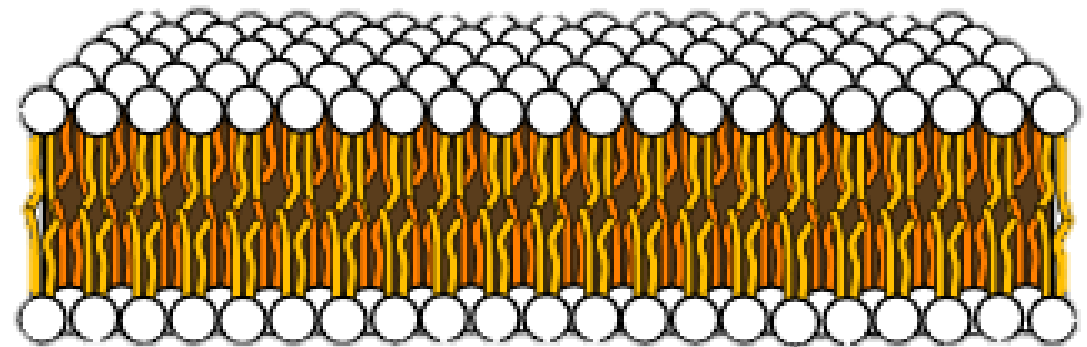- Colors refer to intensity/ density of neutrons scattered

- Has a hydrophobic tail, hydrophilic head
- Form variety of shapes
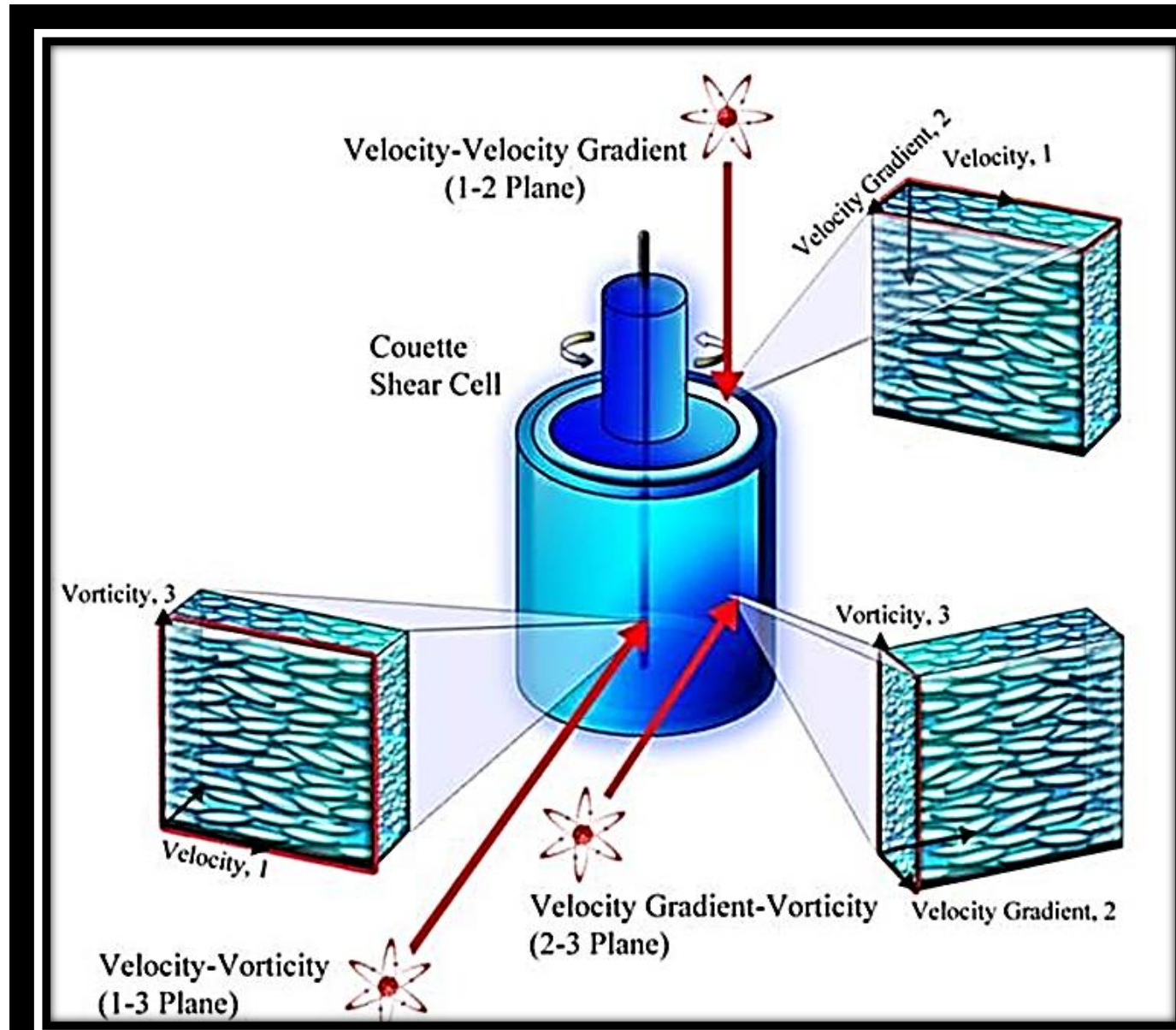- Used in soap, drug delivery, even mayonnaise!

Liposome
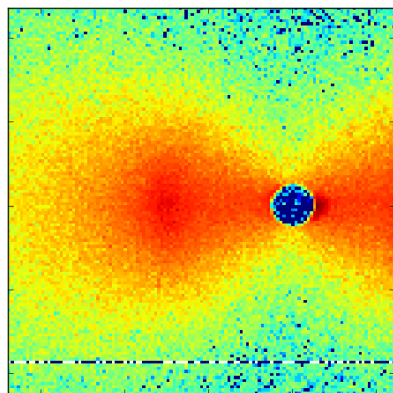
Micelle

Bilayer Sheet/Lamellar

- Device spins the surfactants at different speeds
- Surfactants change shape
- How does the geometry of the material change under different frequencies?
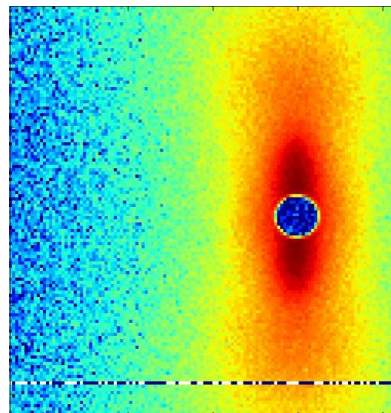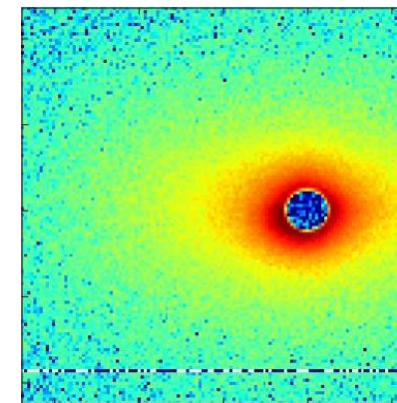
- Fitting allows a better understanding of the data
- Chose a model (such as Core-Shell-Cylinder, Lamellar, Triaxial-Ellipse) to best describe the phase of the surfactant micelles
- Manipulate variables to fit the shape and orientation of the data

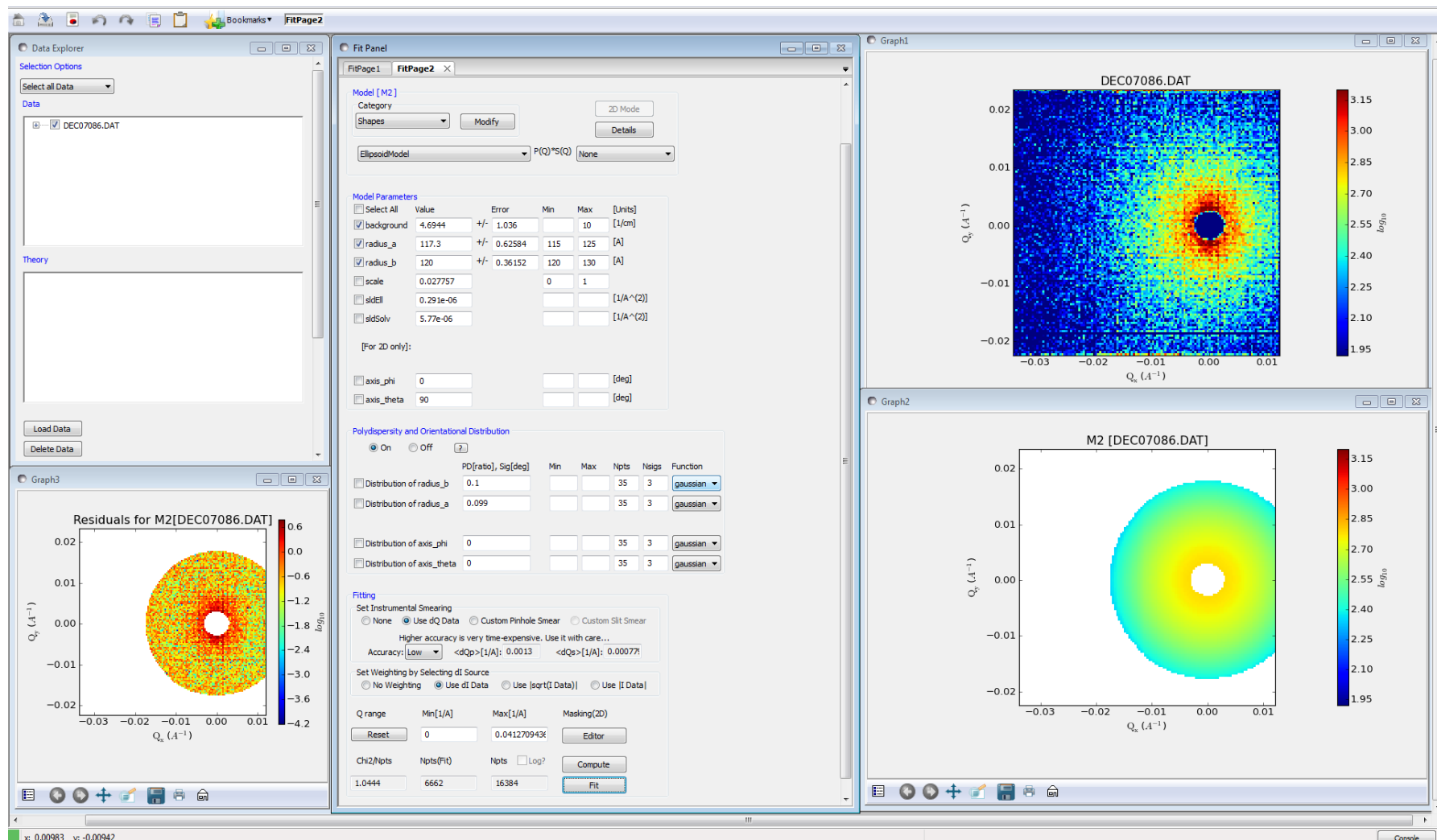➢ These are examples of data collected under various amounts of shear stress



**0 Hz**



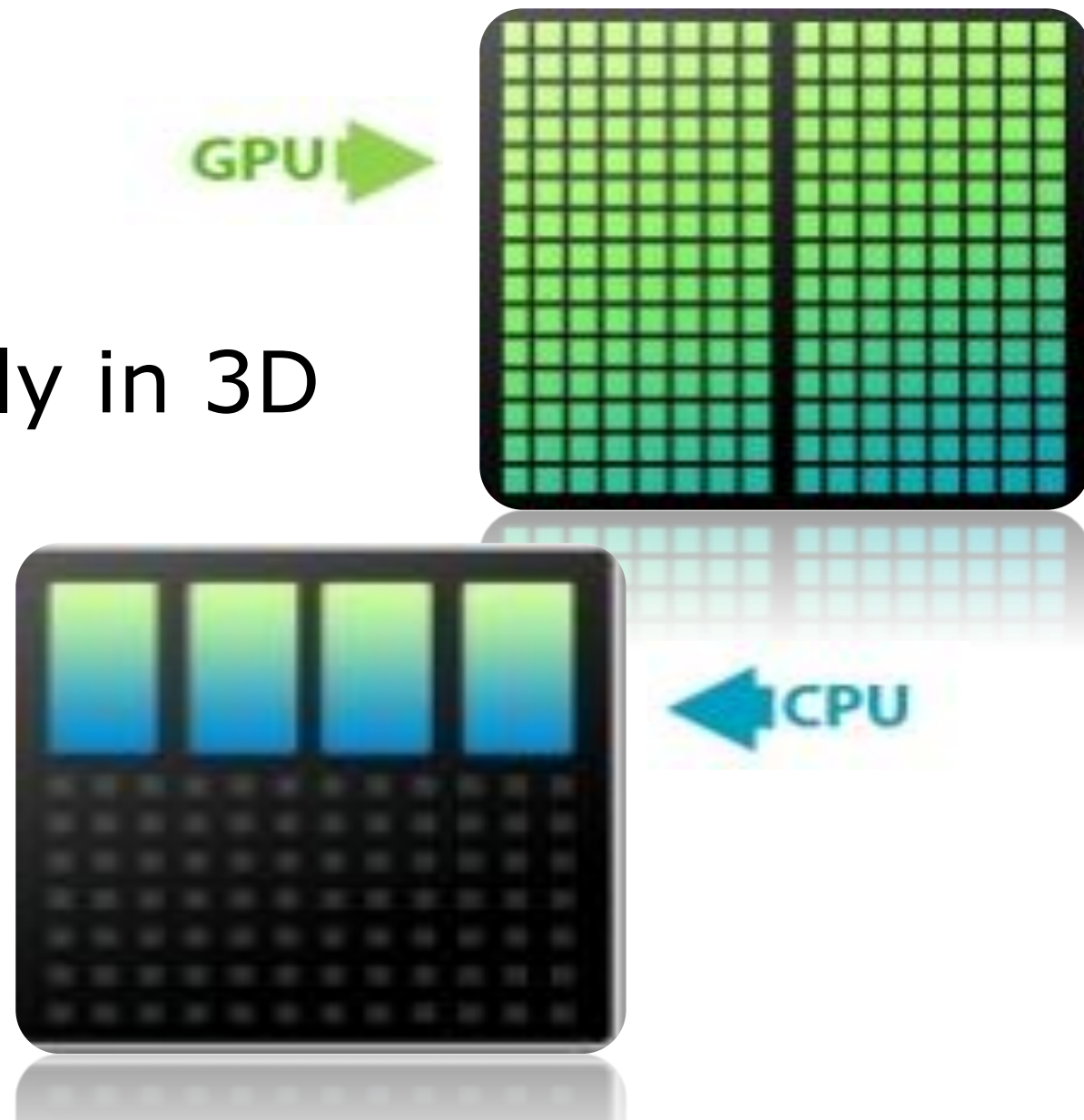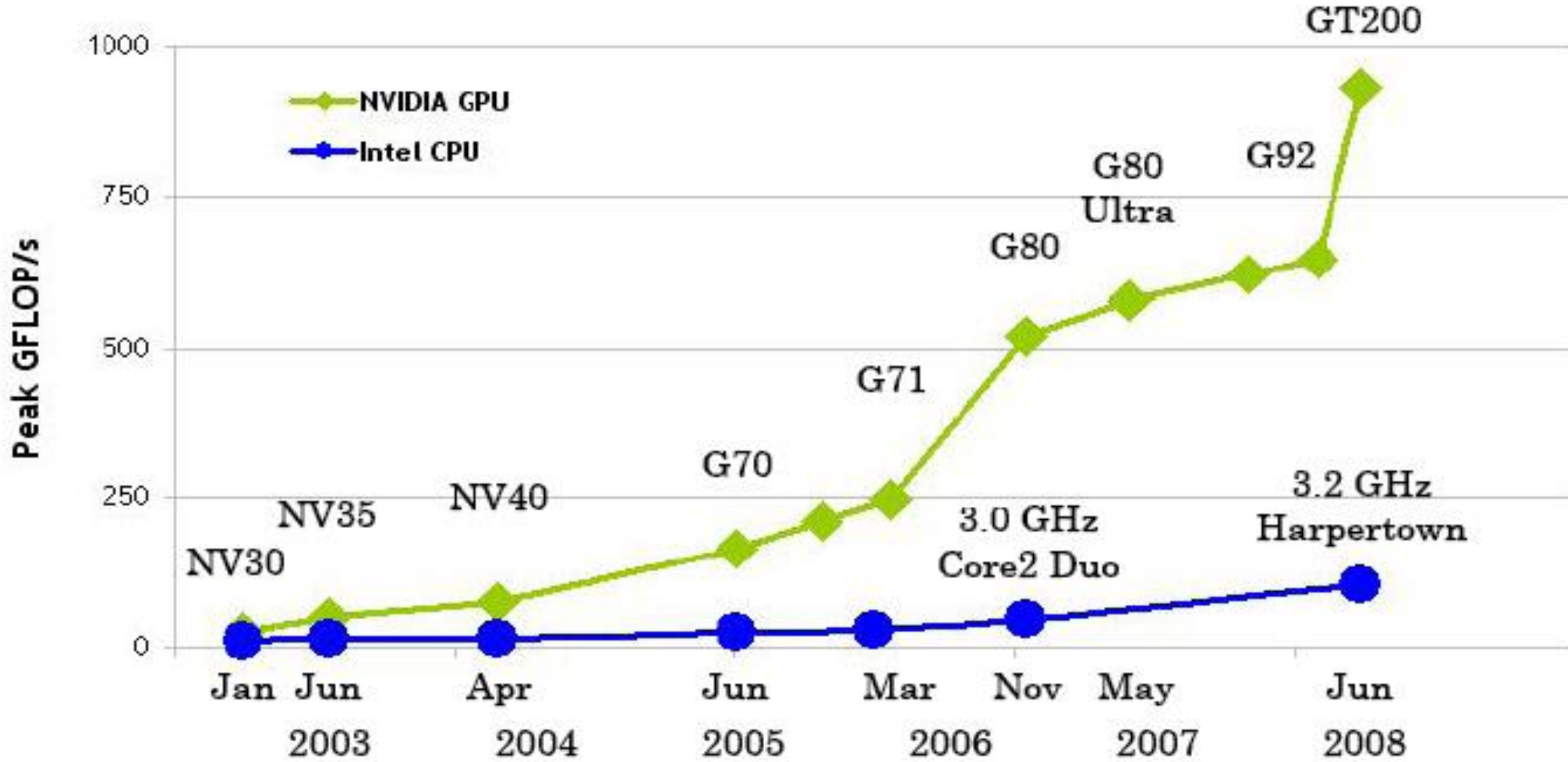**1230 Hz**



**7000 Hz**

# SASVIEW: FITS THE DATA



- Choose model that best describes data

- Issues: Slow, crashes often

# SOLUTION!

- GPU programming: faster
- CPU—found in most PCs
- GPU—previously used solely in 3D gaming
- GPU allows parallel processing, 1000s of threads, 100s more cores than CPU

- GPUs are cheap, fast, and energy efficient

CPU

GPU

| Processor | US$ | Cores | GHz | GFlops | GFlops/W |
|---|---|---|---|---|---|
| Intel i7 x6 | 1000 | 6 | 3.3 | 100 | <1 |
| AMD Phenom II x6 | 300 | 6 | 4.0 | 100 | <1 |
| NVidia GeForce 480 | 500 | 480 | 0.7 | 1400 | 5 |
| ATI Radeon 5970 | 700 | 3200 | 0.7 | 4600 | 15 |

Speed-ups of different projects

- For matrix multiplication: **GPU = 150*CPU MFLOPS** (Mega floating point operations per second)

- **But**, need to tune algorithm and memory transfer for every sized GPU, and for each kernel and program

| | |
|---|---|
| MA Hospital | 300X |
| U Rochester | 160X |
| U Amsterdam | 150X |
| Harvard | 130X |
| U Pennsylvania | 130X |
| Nanyang Tech | 130X |
| U Illinois | 125X |
| Cambridge U | 100X |
| Boise State | 100X |
| Florida U | 100X |

# HOST

The Host (CPU) passes memory buffers, kernels, and queue commands to the device, and receives the result, also in a buffer

# HOST
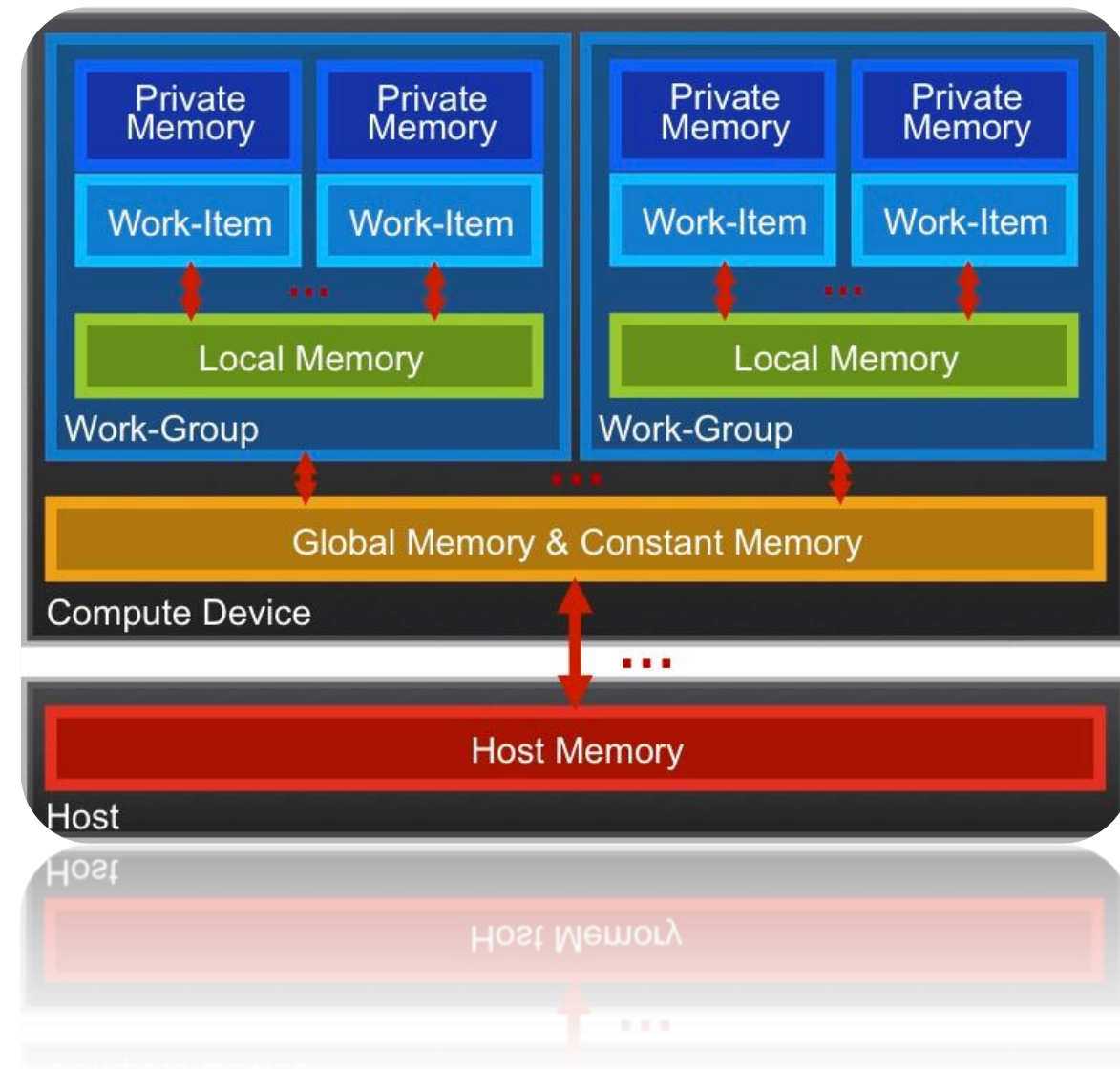
## CONTEXT

The Context holds the GPU(s)—varying shaped & sizes

A Device (GPU) handles the computations.

**Handling of memory transfer**
**global → local → private**
**effect speeds**

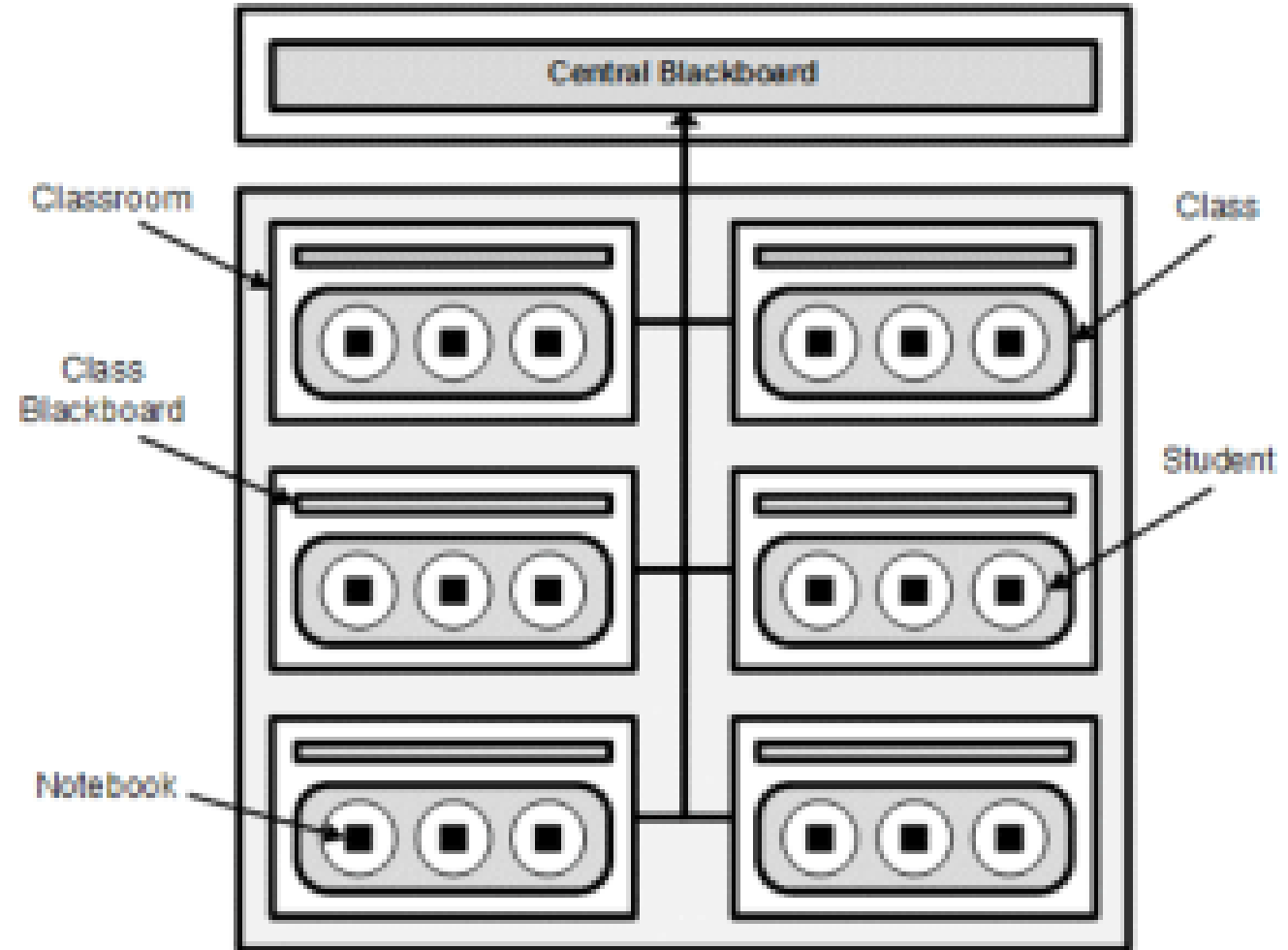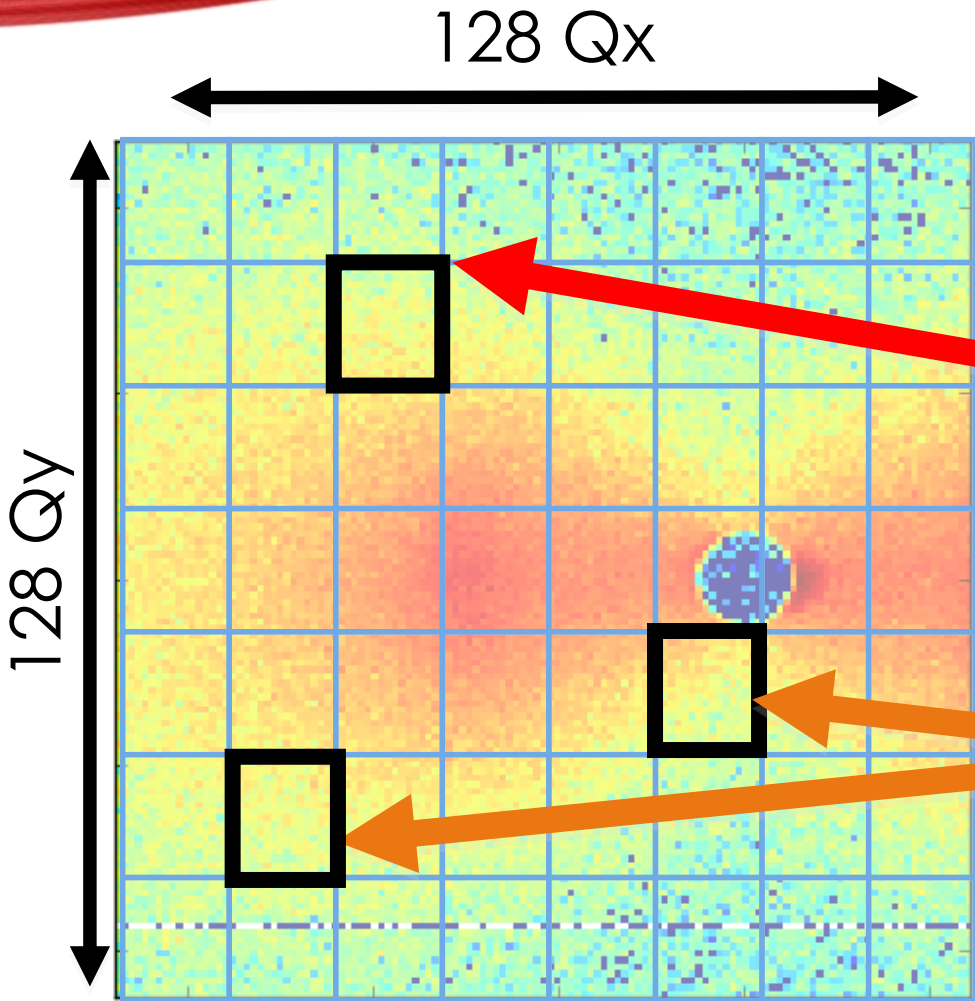A Device has global, local, and private memory, and many work groups that perform calculations in parallel

Central Blackboard: Global Memory

Class Blackboard: Local Memory

Notebook: Private Memory

128 Qx

128 Qy

**A work group,
128x128 performs
computations**

**Work-groups**
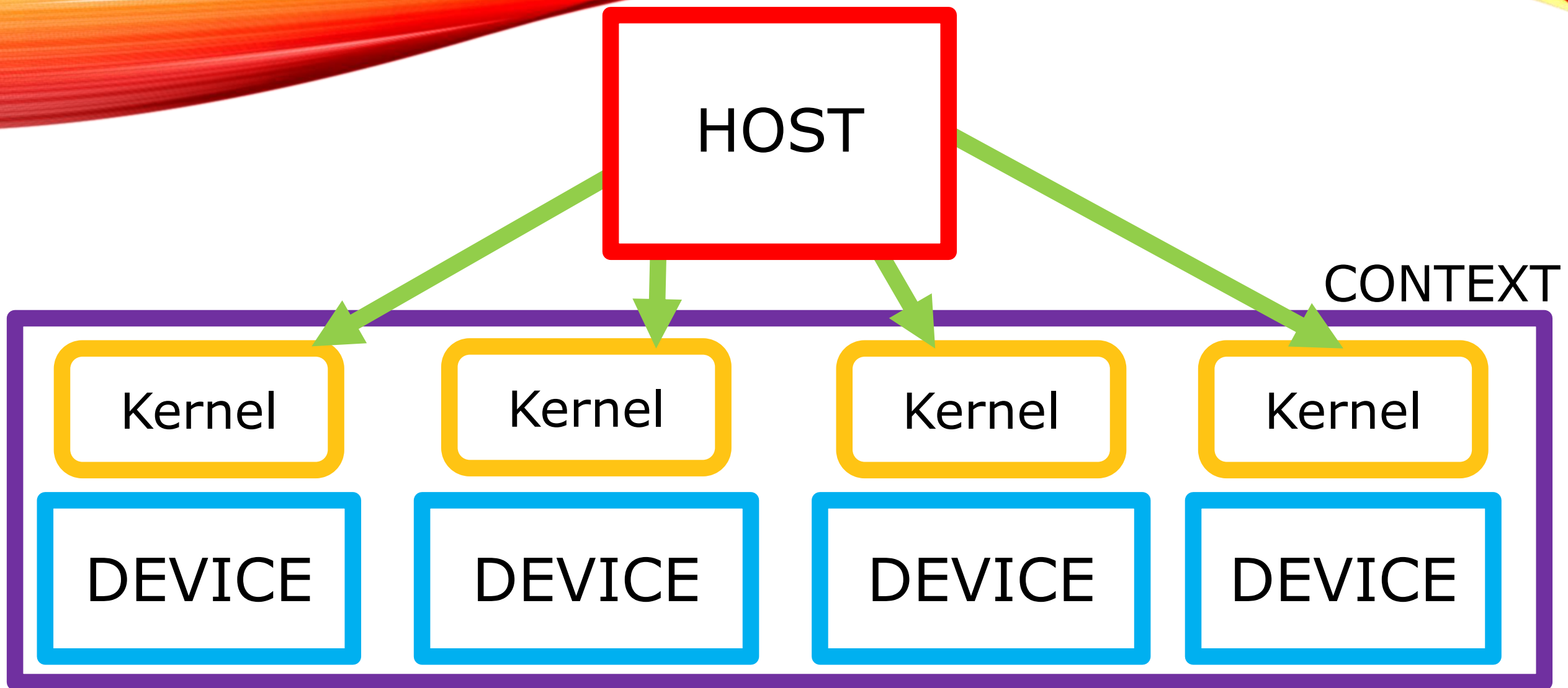**run in parallel,
do not interact**

HOST

CONTEXT

Kernel

Kernel

Kernel

Kernel

DEVICE

DEVICE

DEVICE

DEVICE

The Kernel is the physical code, or program, that are computed by the device

# Our Kernel: calculates scattering intensity, stored in a 2D array

- In our algorithm, work-group takes random Qx &Qy and calculates the scattering density at that point
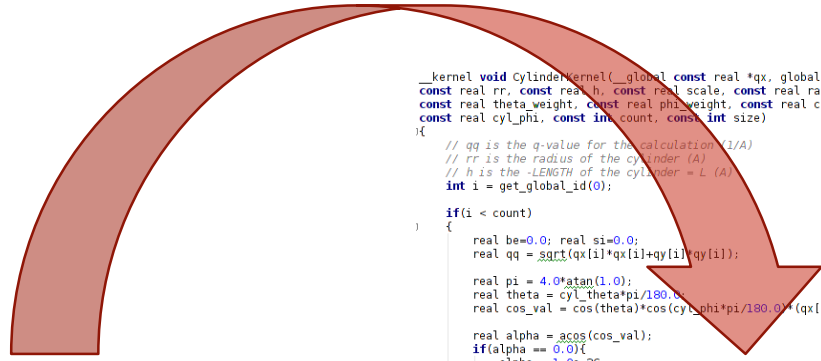- Performs this until every Qx & Qy complete
- Adds results and returns to CPU

Cylinder Model kernel:



Cylinder Model equations:

$$P(q, \alpha) = \frac{scale}{V_s} f^2(q) + bkg$$

$$f(q) = 2(\rho_c - \rho_s)V_c \sin[qL \cos \alpha/2]/[qL \cos \alpha/2]\frac{J_1[qr \sin \alpha]}{[qr \sin \alpha]}$$

$$+ 2(\rho_s - \rho_{solv})V_s \sin[q(L + 2t) \cos \alpha/2]/[q(L + 2t) \cos \alpha/2]\frac{J_1[q(r + t) \sin \alpha]}{[q(r + t) \sin \alpha]}$$

$$V_s = \pi(R + t)^2 \cdot (L + 2t)$$

Loop for polydispersity in CPU
- Size of polydispersity corresponds to width of bell-curve
- Allow a variety of values for a variable (like length)
- For example, high polydispersity in theta gives a larger range of angles
- Also, the more bins, the more accurate the fit

← High polydispersity, but
   low number of bins (5)
→ Lower polydispersity,
   but many bins (40)

Cost: higher bins
means much slower
fit, so need to balance

- In CPU, program the context, device(s), the queue to relay information, and write the buffers for variables, and return values
- Using bumps, loop (again!) over the entire program to fit different variables

```
75
76    model = SasModel(data, GpuCylinder,
77    scale=0.0104,
78    radius=92.5,
79    length=798.3,
80    sldCyl=.29e-6,
81    sldSolv=7.105e-6,
82    background=5,
83    cyl_theta=0,
84    cyl_phi=0,
85    cyl_theta_pd=22.11,
86    cyl_theta_pd_n=20,
87    cyl_theta_pd_nsigma=3,
88    radius_pd=.0084,
89    radius_pd_n=10,
90    radius_pd_nsigma=3,
91    length_pd=0.493,
92    length_pd_n=10,
93    length_pd_nsigma=3,
94    cyl_phi_pd=0,
95    cyl_phi_pd_n=1,
96    cyl_phi_pd_nsigma=3,
97    dtype='f
98
```
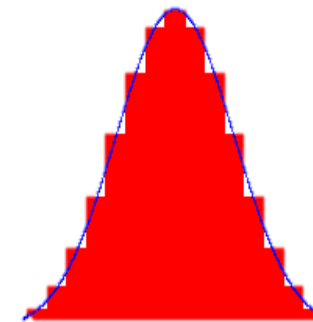
```
99
00
01    # SET THE FITTING PARAMETERS
02    model.radius.range(1, 500)
03    model.length.range(1, 4000)
04    model.cyl_theta.range(-90,100)
05    model.cyl_theta_pd.range(0, 90)
06    model.cyl_theta_pd_n = model.cyl_theta_pd + 5
07    model.radius_pd.range(0, 90)
08    model.length_pd.range(0, 90)
09    model.scale.range(0, 1)
10    model.background.range(0, 100)
11    model.sldCyl.range(0, 1)
12
13
14
```

| Model | Sasview | GPU | Speedup |
|---|---|---|---|
| Cylinder | 3977.7ms | 202.3ms | 19.7X |
| Ellipse | 2953.2ms | 285.5ms | 10.3X |
| Core-Shell-Cylind | 71149.9ms | 4474.7ms | 15.9X |
| Triaxial-Ellipse | 100627ms | 6500.2ms | 15.5X |
| Lamellar | 69.2ms | 6.2ms | 11.2X |

- **Day-long fit to hour-long fit**

- Paul: **50X** faster—cuts out values when the polydispersity weight is low, use local memory

- If 4 GPUs: 4 times faster (200X)

- Allows increased control over simultaneous fitting, multiple-model fitting, and the angular limits of integration in 1D

- Also used models to fit various scattering data

- Here is an example of a fit for a surfactant at 0 Hz
- The left is the data, the middle is the fit, and the right is the residuals of the data

# ACKNOWLEDGEMENTS

❖ Paul Kienzle

❖ Dr. Matt Wasbrough

❖ Aaron West

❖ Yusuf Ameri cainaki

# OPENCL VS CUDA

CUDA:
- Easier to understand; more tutorials online and in books
- CUDA: need to have whole toolchain available

OpenCL:
- Newer, so not much online; less accessible to learn
- Broader range of hardware supported
- simply link the shared library to access